*The state of the art in data compression is arithmetic coding, not the better-known Huffman method. Arithmetic coding gives greater compression, is faster for adaptive models, and clearly separates the model from the channel encoding.*

# ARITHMETIC CODING FOR DATA COMPRESSION

## IAN H. WITTEN, RADFORD M. NEAL, and JOHN G. CLEARY

Arithmetic coding is superior in most respects to the better-known Huffman [10] method. It represents information at least as compactly—sometimes considerably more so. Its performance is optimal without the need for blocking of input data. It encourages a clear separation between the model for representing data and the encoding of information with respect to that model. It accommodates adaptive models easily and is computationally efficient. Yet many authors and practitioners seem unaware of the technique. Indeed there is a widespread belief that Huffman coding cannot be improved upon.

We aim to rectify this situation by presenting an accessible implementation of arithmetic coding and by detailing its performance characteristics. We start by briefly reviewing basic concepts of data compression and introducing the model-based approach that underlies most modern techniques. We then outline the idea of arithmetic coding using a simple example, before presenting programs for both encoding and decoding. In these programs the model occupies a separate module so that different models can easily be used. Next we discuss the construction of fixed and adaptive models and detail the compression efficiency and execution time of the programs, including the effect of different arithmetic word lengths on compression efficiency. Finally, we outline a few applications where arithmetic coding is appropriate.

## DATA COMPRESSION

To many, data compression conjures up an assortment of ad hoc techniques such as conversion of spaces in text to tabs, creation of special codes for common words, or run-length coding of picture data (e.g., see [8]). This contrasts with the more modern model-based paradigm for coding, where, from an *input string* of symbols and a *model*, an *encoded string* is produced that is (usually) a compressed version of the input. The decoder, which must have access to the same model, regenerates the exact input string from the encoded string. Input symbols are drawn from some well-defined set such as the ASCII or binary alphabets; the encoded string is a plain sequence of bits. The model is a way of calculating, in any given context, the distribution of probabilities for the next input symbol. It must be possible for the decoder to produce exactly the same probability distribution in the same context. Compression is achieved by transmitting the more probable symbols in fewer bits than the less probable ones.

For example, the model may assign a predetermined probability to each symbol in the ASCII alphabet. No context is involved. These probabilities can be determined by counting frequencies in representative samples of text to be transmitted. Such a *fixed* model is communicated in advance to both encoder and decoder, after which it is used for many messages.

Alternatively, the probabilities that an *adaptive* model assigns may change as each symbol is transmitted, based on the symbol frequencies seen so far in the message. There is no need for a representative

sample of text, because each message is treated as an independent unit, starting from scratch. The encoder's model changes with each symbol transmitted, and the decoder's changes with each symbol received, in sympathy.

More complex models can provide more accurate probabilistic predictions and hence achieve greater compression. For example, several characters of previous context could condition the next-symbol probability. Such methods have enabled mixed-case English text to be encoded in around 2.2 bits/character with two quite different kinds of model [4, 6]. Techniques that do not separate modeling from coding so distinctly, like that of Ziv and Lempel [23], do not seem to show such great potential for compression, although they may be appropriate when the aim is raw speed rather than compression performance [22].

The effectiveness of any model can be measured by the entropy of the message with respect to it, usually expressed in bits/symbol. Shannon's fundamental theorem of coding states that, given messages randomly generated from a model, it is impossible to encode them into less bits (on average) than the entropy of that model [21].

A message can be coded with respect to a model using either Huffman or arithmetic coding. The former method is frequently advocated as the best possible technique for reducing the encoded data rate. It is not. Given that each symbol in the alphabet must translate into an integral number of bits in the encoding, Huffman coding indeed achieves "minimum redundancy." In other words, it performs optimally if all symbol probabilities are integral powers of ½. But this is not normally the case in practice; indeed, Huffman coding can take up to one extra bit per symbol. The worst case is realized by a source in which one symbol has probability approaching unity. Symbols emanating from such a source convey negligible information on average, but require at least one bit to transmit [7]. Arithmetic coding dispenses with the restriction that each symbol must translate into an integral number of bits, thereby coding more efficiently. It actually achieves the theoretical entropy bound to compression efficiency for any source, including the one just mentioned.

In general, sophisticated models expose the deficiencies of Huffman coding more starkly than simple ones. This is because they more often predict symbols with probabilities close to one, the worst case for Huffman coding. For example, the techniques mentioned above that code English text in 2.2 bits/character both use arithmetic coding as the final step, and performance would be impacted severely

if Huffman coding were substituted. Nevertheless, since our topic is coding and not modeling, the illustrations in this article all employ simple models. Even so, as we shall see, Huffman coding is inferior to arithmetic coding.

The basic concept of arithmetic coding can be traced back to Elias in the early 1960s (see [1, pp. 61–62]). Practical techniques were first introduced by Rissanen [16] and Pasco [15], and developed further by Rissanen [17]. Details of the implementation presented here have not appeared in the literature before; Rubin [20] is closest to our approach. The reader interested in the broader class of arithmetic codes is referred to [18]; a tutorial is available in [13]. Despite these publications, the method is not widely known. A number of recent books and papers on data compression mention it only in passing, or not at all.

## THE IDEA OF ARITHMETIC CODING

In arithmetic coding, a message is represented by an interval of real numbers between 0 and 1. As the message becomes longer, the interval needed to represent it becomes smaller, and the number of bits needed to specify that interval grows. Successive symbols of the message reduce the size of the interval in accordance with the symbol probabilities generated by the model. The more likely symbols reduce the range by less than the unlikely symbols and hence add fewer bits to the message.

Before anything is transmitted, the range for the message is the entire interval [0, 1), denoting the half-open interval $0 \leq x < 1$. As each symbol is processed, the range is narrowed to that portion of it allocated to the symbol. For example, suppose the alphabet is {a, e, i, o, u, !}, and a fixed model is used with probabilities shown in Table I. Imagine trans-

**TABLE I. Example Fixed Model for Alphabet {a, e, i, o, u, !}**

| Symbol | Probability | Range |
|--------|-------------|------------|
| a | .2 | [0, 0.2) |
| e | .3 | [0.2, 0.5) |
| i | .1 | [0.5, 0.6) |
| o | .2 | [0.6, 0.8) |
| u | .1 | [0.8, 0.9) |
| ! | .1 | [0.9, 1.0) |

mitting the message *eaii!*. Initially, both encoder and decoder know that the range is [0, 1). After seeing the first symbol, *e*, the encoder narrows it to [0.2, 0.5), the range the model allocates to this symbol. The second symbol, *a*, will narrow this new range to the first one-fifth of it, since *a* has been

allocated [0, 0.2). This produces [0.2, 0.26), since the previous range was 0.3 units long and one-fifth of that is 0.06. The next symbol, *i*, is allocated [0.5, 0.6), which when applied to [0.2, 0.26) gives the smaller range [0.23, 0.236). Proceeding in this way, the encoded message builds up as follows:

| | | |
|---|---|---|
| Initially | [0, | 1) |
| After seeing *e* | [0.2, | 0.5) |
| *a* | [0.2, | 0.26) |
| *i* | [0.23, | 0.236) |
| *i* | [0.233, | 0.2336) |
| *!* | [0.23354, | 0.2336) |

Figure 1 shows another representation of the encoding process. The vertical bars with ticks represent the symbol probabilities stipulated by the model. After the first symbol has been processed, the model is scaled into the range [0.2, 0.5), as shown in
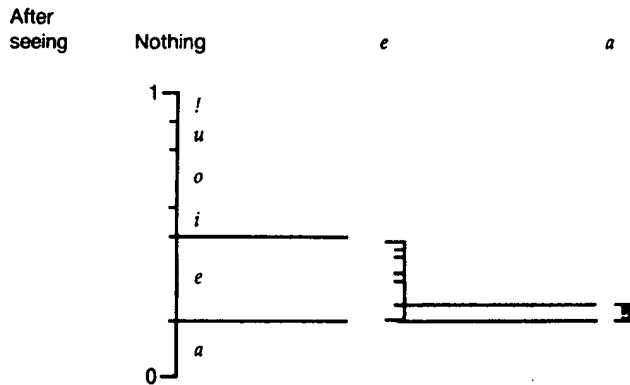
Figure 1a. The second symbol scales it again into the range [0.2, 0.26). But the picture cannot be continued in this way without a magnifying glass! Consequently, Figure 1b shows the ranges expanded to full height at every stage and marked with a scale that gives the endpoints as numbers.

Suppose all the decoder knows about the message is the final range, [0.23354, 0.2336). It can immediately deduce that the first character was *e*, since the range lies entirely within the space the model of Table I allocates for *e*. Now it can simulate the operation of the *e*ncoder:

| | |
|---|---|
| Initially | [0, 1) |
| After seeing *e* | [0.2, 0.5) |

This makes it clear that the second character is *a*, since this will produce the range

| | |
|---|---|
| After seeing *a* | [0.2, 0.26), |

which entirely encloses the given range [0.23354, 0.2336). Proceeding like this, the decoder can identify the whole message.

It is not really necessary for the decoder to know both ends of the range produced by the encoder. Instead, a single number within the range—for example, 0.23355—will suffice. (Other numbers, like 0.23354, 0.23357, or even 0.23354321, would do just as well.) However, the decoder will face the problem of detecting the end of the message, to determine when to stop decoding. After all, the single number 0.0 could represent any of *a*, *aa*, *aaa*, *aaaa*, . . . . To resolve the ambiguity, we ensure that each message ends with a special EOF symbol known to both encoder and decoder. For the alphabet of Table I, *!* will be used to terminate messages, and only to termi-



**After seeing Nothing** *e* *a*

**FIGURE 1a. Representation of the Arithmetic Coding Process**
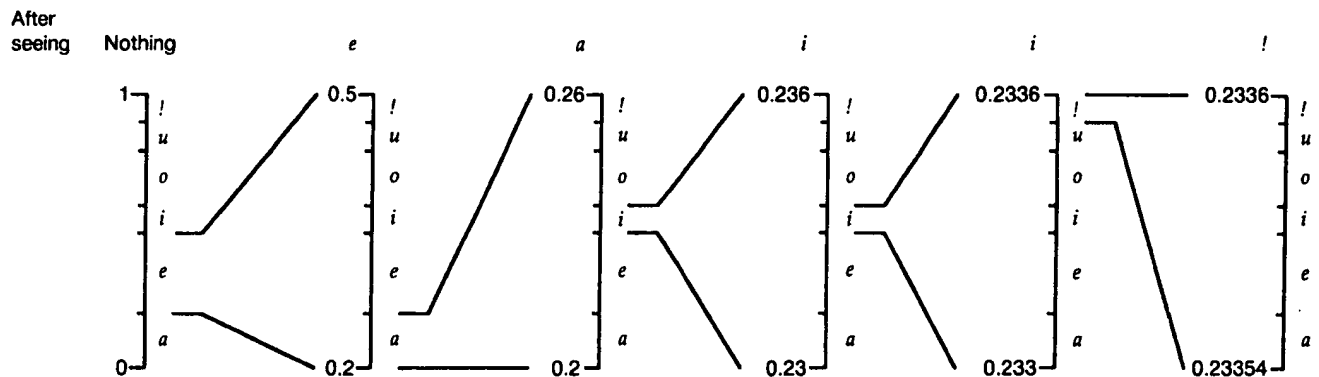


**After seeing Nothing** *e* *a* *i* *i* *!*

**FIGURE 1b. Representation of the Arithmetic Coding Process with the Interval Scaled Up at Each Stage**

```
/* ARITHMETIC ENCODING ALGORITHM. */

/* Call encode_symbol repeatedly for each symbol in the message.          */
/* Ensure that a distinguished "terminator" symbol is encoded last, then  */
/* transmit any value in the range [low, high).                           */

encode_symbol(symbol, cum_freq)
    range = high - low
    high  = low + range*cum_freq[symbol-1]
    low   = low + range*cum_freq[symbol]


/* ARITHMETIC DECODING ALGORITHM. */

/* "Value" is the number that has been received.                         */
/* Continue calling decode_symbol until the terminator symbol is returned. */

decode_symbol(cum_freq)
    find symbol such that
        cum_freq[symbol] <= (value-low)/(high-low) < cum_freq[symbol-1]
                         /* This ensures that value lies within the new */
                         /* [low, high) range that will be calculated by */
                         /* the following lines of code.                 */

    range = high - low
    high  = low + range*cum_freq[symbol-1]
    low   = low + range*cum_freq[symbol]
    return symbol
```

**FIGURE 2. Pseudocode for the Encoding and Decoding Procedures**

nate messages. When the decoder sees this symbol, it stops decoding.

Relative to the fixed model of Table I, the entropy of the five-symbol message *eaii!* is

$$-\log 0.3 - \log 0.2 - \log 0.1 - \log 0.1 - \log 0.1$$

$$= -\log 0.00006 \approx 4.22$$

(using base 10, since the above encoding was performed in decimal). This explains why it takes five decimal digits to encode the message. In fact, the size of the final range is $0.2336 - 0.23354 = 0.00006$, and the entropy is the negative logarithm of this figure. Of course, we normally work in binary, transmitting binary digits and measuring entropy in bits.

Five decimal digits seems a lot to encode a message comprising four vowels! It is perhaps unfortunate that our example ended up by expanding rather than compressing. Needless to say, however, different models will give different entropies. The best single-character model of the message *eaii!* is the set of symbol frequencies {*e*(0.2), *a*(0.2), *i*(0.4), *!*(0.2)}, which gives an entropy of 2.89 decimal digits. Using this model the encoding would be only three digits long. Moreover, as noted earlier, more sophisticated models give much better performance in general.

## A PROGRAM FOR ARITHMETIC CODING

Figure 2 shows a pseudocode fragment that summarizes the encoding and decoding procedures developed in the last section. Symbols are numbered, 1, 2, 3, .... The frequency range for the $i$th symbol is from $cum\_freq[i]$ to $cum\_freq[i-1]$. As $i$ decreases, $cum\_freq[i]$ increases, and $cum\_freq[0] = 1$. (The reason for this "backwards" convention is that $cum\_freq[0]$ will later contain a normalizing factor, and it will be convenient to have it begin the array.) The "current interval" is $[low, high)$, and for both encoding and decoding, this should be initialized to $[0, 1)$.

Unfortunately, Figure 2 is overly simplistic. In practice, there are several factors that complicate both encoding and decoding:

*Incremental transmission and reception.* The encode algorithm as described does not transmit anything until the entire message has been encoded; neither does the decode algorithm begin decoding until it has received the complete transmission. In most applications an incremental mode of operation is necessary.

*The desire to use integer arithmetic.* The precision required to represent the $[low, high)$ interval grows with the length of the message. Incremental operation will help overcome this, but the potential for

overflow and underflow must still be examined carefully.

*Representing the model so that it can be consulted efficiently.* The representation used for the model should minimize the time required for the decode algorithm to identify the next symbol. Moreover, an adaptive model should be organized to minimize the time-consuming task of maintaining cumulative frequencies.

Figure 3 shows working code, in C, for arithmetic encoding and decoding. It is considerably more detailed than the bare-bones sketch of Figure 2! Implementations of two different models are given in Figure 4; the Figure 3 code can use either one.

The remainder of this section examines the code of Figure 3 more closely, and includes a proof that decoding is still correct in the integer implementation and a review of constraints on word lengths in the program.

```
arithmetic_coding.h
```

```
1   /* DECLARATIONS USED FOR ARITHMETIC ENCODING AND DECODING */
2
3
4   /* SIZE OF ARITHMETIC CODE VALUES. */
5
6   #define Code_value_bits 16          /* Number of bits in a code value   */
7   typedef long code_value;            /* Type of an arithmetic code value */
8
9   #define Top_value (((long)1<<Code_value_bits)-1)      /* Largest code value */
10
11
12  /* HALF AND QUARTER POINTS IN THE CODE VALUE RANGE. */
13
14  #define First_qtr (Top_value/4+1)   /* Point after first quarter   */
15  #define Half      (2*First_qtr)     /* Point after first half      */
16  #define Third_qtr (3*First_qtr)     /* Point after third quarter   */
```

```
model.h
```

```
17  /* INTERFACE TO THE MODEL. */
18
19
20  /* THE SET OF SYMBOLS THAT MAY BE ENCODED. */
21
22  #define No_of_chars 256             /* Number of character symbols   */
23  #define EOF_symbol (No_of_chars+1)  /* Index of EOF symbol           */
24
25  #define No_of_symbols (No_of_chars+1)  /* Total number of symbols    */
26
27
28  /* TRANSLATION TABLES BETWEEN CHARACTERS AND SYMBOL INDEXES. */
29
30  int char_to_index[No_of_chars];          /* To index from character   */
31  unsigned char index_to_char[No_of_symbols+1]; /* To character from index */
32
33
34  /* CUMULATIVE FREQUENCY TABLE. */
35
36  #define Max_frequency 16383         /* Maximum allowed frequency count */
37                                      /*    2^14 - 1                     */
38  int cum_freq[No_of_symbols+1];      /* Cumulative symbol frequencies   */
```

FIGURE 3.   C Implementation of Arithmetic Encoding and Decoding

encode.c

```
39   /* MAIN PROGRAM FOR ENCODING. */
40
41   #include <stdio.h>
42   #include "model.h"
43
44   main()
45   {    start_model();                           /* Set up other modules.      */
46        start_outputing_bits();
47        start_encoding();
48        for (;;) {                               /* Loop through characters. */
49            int ch; int symbol;
50            ch = getc(stdin);                    /* Read the next character. */
51            if (ch==EOF) break;                  /* Exit loop on end-of-file.*/
52            symbol = char_to_index[ch];          /* Translate to an index.     */
53            encode_symbol(symbol,cum_freq);      /* Encode that symbol.       */
54            update_model(symbol);                /* Update the model.         */
55        }
56        encode_symbol(EOF_symbol,cum_freq);      /* Encode the EOF symbol.    */
57        done_encoding();                         /* Send the last few bits.   */
58        done_outputing_bits();
59        exit(0);
60   }
```

arithmetic_encode.c

```
61   /* ARITHMETIC ENCODING ALGORITHM. */
62
63   #include "arithmetic_coding.h"
64
65   static void bit_plus_follow();   /* Routine that follows            */
66
67
68   /* CURRENT STATE OF THE ENCODING. */
69
70   static code_value low, high;     /* Ends of the current code region       */
71   static long bits_to_follow;      /* Number of opposite bits to output after */
72                                    /* the next bit.                    */
73
74
75   /* START ENCODING A STREAM OF SYMBOLS. */
76
77   start_encoding()
78   {    low = 0;                                /* Full code range.         */
79        high = Top_value;
80        bits_to_follow = 0;                     /* No bits to follow next.  */
81   }
82
83
84   /* ENCODE A SYMBOL. */
85
86   encode_symbol(symbol,cum_freq)
87        int symbol;                /* Symbol to encode                 */
88        int cum_freq[];            /* Cumulative symbol frequencies    */
89   {    long range;                /* Size of the current code region  */
90        range = (long)(high-low)+1;
91        high = low +                             /* Narrow the code region   */
92          (range*cum_freq[symbol-1])/cum_freq[0]-1; /* to that allotted to this */
93        low = low +                              /* symbol.                  */
94          (range*cum_freq[symbol])/cum_freq[0];
```

FIGURE 3.   C Implementation of Arithmetic Encoding and Decoding (*continued*)

```
95      for (;;) {                                /* Loop to output bits.       */
96          if (high<Half) {
97              bit_plus_follow(0);               /* Output 0 if in low half. */
98          }
99          else if (low>=Half) {                 /* Output 1 if in high half.*/
100             bit_plus_follow(1);
101             low -= Half;
102             high -= Half;                     /* Subtract offset to top.   */
103         }
104         else if (low>=First_qtr               /* Output an opposite bit    */
105                  && high<Third_qtr) {         /* later if in middle half. */
106             bits_to_follow += 1;
107             low -= First_qtr;                 /* Subtract offset to middle*/
108             high -= First_qtr;
109         }
110         else break;                           /* Otherwise exit loop.      */
111         low = 2*low;
112         high = 2*high+1;                      /* Scale up code range.      */
113     }
114  }
115
116
117  /* FINISH ENCODING THE STREAM. */
118
119  done_encoding()
120  {   bits_to_follow += 1;                      /* Output two bits that      */
121      if (low<First_qtr) bit_plus_follow(0);   /* select the quarter that   */
122      else bit_plus_follow(1);                 /* the current code range    */
123  }                                            /* contains.                 */
124
125
126  /* OUTPUT BITS PLUS FOLLOWING OPPOSITE BITS. */
127
128  static void bit_plus_follow(bit)
129      int bit;
130  {   output_bit(bit);                          /* Output the bit.           */
131      while (bits_to_follow>0) {
132          output_bit(!bit);                     /* Output bits_to_follow     */
133          bits_to_follow -= 1;                  /* opposite bits. Set        */
134      }                                         /* bits_to_follow to zero.   */
135  }
```

decode.c

```
136  /* MAIN PROGRAM FOR DECODING. */
137
138  #include <stdio.h>
139  #include "model.h"
140
141  main()
142  {   start_model();                            /* Set up other modules.     */
143      start_inputing_bits();
144      start_decoding();
145      for (;;) {                                /* Loop through characters. */
146          int ch; int symbol;
147          symbol = decode_symbol(cum_freq);     /* Decode next symbol.       */
148          if (symbol==EOF_symbol) break;        /* Exit loop if EOF symbol. */
149          ch = index_to_char[symbol];           /* Translate to a character.*/
150          putc(ch,stdout);                      /* Write that character.     */
151          update_model(symbol);                 /* Update the model.         */
152      }
153      exit(0);
154  }
```

**FIGURE 3.   C Implementation of Arithmetic Encoding and Decoding (*continued*)**

arithmetic_decode.c

```
155    /* ARITHMETIC DECODING ALGORITHM. */
156
157    #include "arithmetic_coding.h"
158
159
160    /* CURRENT STATE OF THE DECODING. */
161
162    static code_value value;           /* Currently-seen code value       */
163    static code_value low, high;       /* Ends of current code region     */
164
165
166    /* START DECODING A STREAM OF SYMBOLS. */
167
168    start_decoding()
169    {   int i;
170        value = 0;                                  /* Input bits to fill the  */
171        for (i = 1; i<=Code_value_bits; i++) {      /* code value.             */
172            value = 2*value+input_bit();
173        }
174        low = 0;                                    /* Full code range.        */
175        high = Top_value;
176    }
177
178
179    /* DECODE THE NEXT SYMBOL. */
180
181    int decode_symbol(cum_freq)
182        int cum_freq[];                /* Cumulative symbol frequencies   */
183    {   long range;                    /* Size of current code region     */
184        int cum;                       /* Cumulative frequency calculated  */
185        int symbol;                    /* Symbol decoded                  */
186        range = (long)(high-low)+1;
187        cum =                                       /* Find cum freq for value. */
188            (((long)(value-low)+1)*cum_freq[0]-1)/range;
189        for (symbol = 1; cum_freq[symbol]>cum; symbol++) ; /* Then find symbol. */
190        high = low +                                /* Narrow the code region  */
191            (range*cum_freq[symbol-1])/cum_freq[0]-1; /* to that allotted to this */
192        low = low +                                 /* symbol.                 */
193            (range*cum_freq[symbol])/cum_freq[0];
194        for (;;) {                                  /* Loop to get rid of bits. */
195            if (high<Half) {
196                /* nothing */                       /* Expand low half.        */
197            }
198            else if (low>=Half) {                   /* Expand high half.       */
199                value -= Half;
200                low -= Half;                        /* Subtract offset to top. */
201                high -= Half;
202            }
203            else if (low>=First_qtr                 /* Expand middle half.     */
204                    && high<Third_qtr) {
205                value -= First_qtr;
206                low -= First_qtr;                   /* Subtract offset to middle*/
207                high -= First_qtr;
208            }
209            else break;                             /* Otherwise exit loop.    */
210            low = 2*low;
211            high = 2*high+1;                        /* Scale up code range.    */
212            value = 2*value+input_bit();            /* Move in next input bit. */
213        }
214        return symbol;
215    }
```

**FIGURE 3.** C Implementation of Arithmetic Encoding and Decoding (*continued*)

bit_input.c

```
216    /* BIT INPUT ROUTINES. */
217
218    #include <stdio.h>
219    #include "arithmetic_coding.h"
220
221
222    /* THE BIT BUFFER. */
223
224    static int buffer;              /* Bits waiting to be input            */
225    static int bits_to_go;          /* Number of bits still in buffer      */
226    static int garbage_bits;        /* Number of bits past end-of-file     */
227
228
229    /* INITIALIZE BIT INPUT. */
230
231    start_inputing_bits()
232    {   bits_to_go = 0;                            /* Buffer starts out with   */
233        garbage_bits = 0;                          /* no bits in it.           */
234    }
235
236
237    /* INPUT A BIT. */
238
239    int input_bit()
240    {   int t;
241        if (bits_to_go==0) {                       /* Read the next byte if no */
242            buffer = getc(stdin);                  /* bits are left in buffer. */
243            if (buffer==EOF) {
244                garbage_bits += 1;                              /* Return arbitrary bits*/
245                if (garbage_bits>Code_value_bits-2) {  /* after eof, but check */
246                    fprintf(stderr,"Bad input file\n"); /* for too many such.   */
247                    exit(-1);
248                }
249            }
250            bits_to_go = 8;
251        }
252        t = buffer&1;                              /* Return the next bit from */
253        buffer >>= 1;                              /* the bottom of the byte.  */
254        bits_to_go -= 1;
255        return t;
256    }
```

FIGURE 3.   C Implementation of Arithmetic Encoding and Decoding (*continued*)

bit_output.c

---

```
257    /* BIT OUTPUT ROUTINES. */
258
259    #include <stdio.h>
260
261
262    /* THE BIT BUFFER. */
263
264    static int buffer;              /* Bits buffered for output              */
265    static int bits_to_go;         /* Number of bits free in buffer         */
266
267
268    /* INITIALIZE FOR BIT OUTPUT. */
269
270    start_outputing_bits()
271    {   buffer = 0;                             /* Buffer is empty to start */
272        bits_to_go= 8;                          /* with.                    */
273    }
274
275
276    /* OUTPUT A BIT. */
277
278    output_bit(bit)
279        int bit;
280    {   buffer >>= 1;                           /* Put bit in top of buffer.*/
281        if (bit) buffer |= 0x80;
282        bits_to_go -= 1;
283        if (bits_to_go==0) {                    /* Output buffer if it is   */
284            putc(buffer,stdout);                /* now full.                */
285            bits_to_go = 8;
286        }
287    }
288
289
290    /* FLUSH OUT THE LAST BITS. */
291
292    done_outputing_bits()
293    {   putc(buffer>>bits_to_go,stdout);
294    }
```

**FIGURE 3.    C Implementation of Arithmetic Encoding and Decoding (*continued*)**

fixed_model.c

```
1    /* THE FIXED SOURCE MODEL */
2
3    #include "model.h"
4
5    int freq[No_of_symbols+1] = {
6        0,
7        1,    1,    1,    1,    1,    1,    1,    1,    1,    1, 124,    1,    1,    1,    1,    1,
8        1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,
9
10   /*          !     "     #     $     %     &     '     (     )     *     +     ,     -     .     /  */
11   1236,    1,   21,    9,    3,    1,   25,   15,    2,    2,    2,    1,   79,   19,   60,    1,
12
13   /*  0     1     2     3     4     5     6     7     8     9     :     ;     <     =     >     ?  */
14     15,   15,    8,    5,    4,    7,    5,    4,    4,    6,    3,    2,    1,    1,    1,    1,
15
16   /*  @     A     B     C     D     E     F     G     H     I     J     K     L     M     N     O  */
17      1,   24,   15,   22,   12,   15,   10,    9,   16,   16,    8,    6,   12,   23,   13,   11,
18
19   /*  P     Q     R     S     T     U     V     W     X     Y     Z     [     \     ]     ^     _  */
20     14,    1,   14,   28,   29,    6,    3,   11,    1,    3,    1,    1,    1,    1,    1,    3,
21
22   /*  `     a     b     c     d     e     f     g     h     i     j     k     l     m     n     o  */
23      1,  491,   85,  173,  232,  744,  127,  110,  293,  418,    6,   39,  250,  139,  429,  446,
24
25   /*  p     q     r     s     t     u     v     w     x     y     z     {     |     }     ~        */
26    111,    5,  388,  375,  531,  152,   57,   97,   12,  101,    5,    2,    1,    2,    3,    1,
27
28      1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,
29      1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,
30      1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,
31      1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,
32      1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,
33      1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,
34      1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,
35      1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,    1,
36      1
37   };
38
39
40   /* INITIALIZE THE MODEL. */
41
42   start_model()
43   {   int i;
44       for (i = 0; i<No_of_chars; i++) {           /* Set up tables that        */
45           char_to_index[i] = i+1;                  /* translate between symbol  */
46           index_to_char[i+1] = i;                  /* indexes and characters.   */
47       }
48       cum_freq[No_of_symbols] = 0;
49       for (i = No_of_symbols; i>0; i--) {          /* Set up cumulative         */
50           cum_freq[i-1] = cum_freq[i] + freq[i];   /* frequency counts.         */
51       }
52       if (cum_freq[0] > Max_frequency) abort();    /* Check counts within limit*/
53   }
54
55
56   /* UPDATE THE MODEL TO ACCOUNT FOR A NEW SYMBOL. */
57
58   update_model(symbol)
59       int symbol;
60   {                                               /* Do nothing. */
61   }
```

**FIGURE 4.  Fixed and Adaptive Models for Use with Figure 3**

adaptive_model.c

```
1    /* THE ADAPTIVE SOURCE MODEL */
2
3    #include "model.h"
4
5    int freq[No_of_symbols+1];        /* Symbol frequencies              */
6
7
8    /* INITIALIZE THE MODEL. */
9
10   start_model()
11   {   int i;
12       for (i = 0; i<No_of_chars; i++) {        /* Set up tables that     */
13           char_to_index[i] = i+1;              /* translate between symbol */
14           index_to_char[i+1] = i;              /* indexes and characters. */
15       }
16       for (i = 0; i<=No_of_symbols; i++) {     /* Set up initial frequency */
17           freq[i] = 1;                         /* counts to be one for all */
18           cum_freq[i] = No_of_symbols-i;       /* symbols.                */
19       }
20       freq[0] = 0;                             /* Freq[0] must not be the */
21   }                                            /* same as freq[1].        */
22
23
24   /* UPDATE THE MODEL TO ACCOUNT FOR A NEW SYMBOL. */
25
26   update_model(symbol)
27       int symbol;                    /* Index of new symbol            */
28   {   int i;                         /* New index for symbol           */
29       if (cum_freq[0]==Max_frequency) {        /* See if frequency counts */
30           int cum;                             /* are at their maximum.  */
31           cum = 0;
32           for (i = No_of_symbols; i>=0; i--) { /* If so, halve all the   */
33               freq[i] = (freq[i]+1)/2;         /* counts (keeping them   */
34               cum_freq[i] = cum;               /* non-zero).             */
35               cum += freq[i];
36           }
37       }
38       for (i = symbol; freq[i]==freq[i-1]; i--) ;  /* Find symbol's new index. */
39       if (i<symbol) {
40           int ch_i, ch_symbol;
41           ch_i = index_to_char[i];             /* Update the translation */
42           ch_symbol = index_to_char[symbol];   /* tables if the symbol has */
43           index_to_char[i] = ch_symbol;        /* moved.                 */
44           index_to_char[symbol] = ch_i;
45           char_to_index[ch_i] = symbol;
46           char_to_index[ch_symbol] = i;
47       }
48       freq[i] += 1;                            /* Increment the frequency */
49       while (i>0) {                            /* count for the symbol and */
50           i -= 1;                              /* update the cumulative  */
51           cum_freq[i] += 1;                    /* frequencies.           */
52       }
53   }
```

**FIGURE 4. Fixed and Adaptive Models for Use with Figure 3 (*continued*)**

### Representing the Model

Implementations of models are discussed in the next section; here we are concerned only with the interface to the model (lines 20–38). In C, a byte is represented as an integer between 0 and 255 (a *char*). Internally, we represent a byte as an integer between 1 and 257 inclusive (an *index*), EOF being treated as a 257th symbol. It is advantageous to sort the model into frequency order, so as to minimize the number of executions of the decoding loop (line 189). To permit such reordering, the *char/index* translation is implemented as a pair of tables, *index_to_char*[ ] and *char_to_index*[ ]. In one of our models, these tables simply form the *index* by adding 1 to the *char*, but another implements a more complex translation that assigns small indexes to frequently used symbols.

The probabilities in the model are represented as integer frequency counts, and cumulative counts are stored in the array *cum_freq*[ ]. As previously, this array is "backwards," and the total frequency count, which is used to normalize all frequencies, appears in *cum_freq*[0]. Cumulative counts must not exceed a predetermined maximum, *Max_frequency*, and the model implementation must prevent overflow by scaling appropriately. It must also ensure that neighboring values in the *cum_freq*[ ] array differ by at least 1; otherwise the affected symbol cannot be transmitted.

### Incremental Transmission and Reception

Unlike Figure 2 the program in Figure 3 represents *low* and *high* as integers. A special data type, *code_value*, is defined for these quantities, together with some useful constants: *Top_value*, representing the largest possible *code_value*, and *First_qtr*, *Half*, and *Third_qtr*, representing parts of the range (lines 6–16). Whereas in Figure 2 the current interval is represented by [*low, high*), in Figure 3 it is [*low, high*]; that is, the range now includes the value of *high*. Actually, it is more accurate (though more confusing) to say that, in the program in Figure 3, the interval represented is [*low, high* + 0.11111 ⋯). This is because when the bounds are scaled up to increase the precision, zeros are shifted into the low-order bits of *low*, but ones are shifted into *high*. Although it is possible to write the program to use a different convention, this one has some advantages in simplifying the code.

As the code range narrows, the top bits of *low* and *high* become the same. Any bits that are the same can be transmitted immediately, since they cannot be affected by future narrowing. For encoding, since we know that *low* ≤ *high*, this requires code like

```
for (;;) {
  if (high < Half) {
    output_bit(0);
    low  = 2*low;
    high = 2*high+1;
  }
  else if (low ≥ Half) {
    output_bit(1);
    low = 2*(low-Half);
    high = 2*(high-Half)+1;
  }
  else break;
}
```

which ensures that, upon completion, *low* < *Half* ≤ *high*. This can be found in lines 95–113 of *encode_symbol*( ), although there are some extra complications caused by underflow possibilities (see the next subsection). Care is taken to shift ones in at the bottom when *high* is scaled, as noted above.

Incremental reception is done using a number called *value* as in Figure 2, in which processed bits flow out the top (high-significance) end and newly received ones flow in the bottom. Initially, *start_decoding*( ) (lines 168–176) fills *value* with received bits. Once *decode_symbol*( ) has identified the next input symbol, it shifts out now-useless high-order bits that are the same in *low* and *high*, shifting *value* by the same amount (and replacing lost bits by fresh input bits at the bottom end):

```
for (;;) {
  if (high < Half) {
    value = 2*value+input_bit( );
    low  = 2*low;
    high = 2*high+1;
  }
  else if (low > Half) {
    value = 2*(value-Half)+input_bit( );
    low  = 2*(low-Half);
    high = 2*(high-Half)+1;
  }
  else break;
}
```

(see lines 194–213, again complicated by precautions against underflow, as discussed below).

### Proof of Decoding Correctness

At this point it is worth checking that identification of the next symbol by *decode_symbol*( ) works properly. Recall from Figure 2 that *decode_symbol*( ) must use *value* to find the symbol that, when encoded, reduces the range to one that still includes *value*. Lines 186–188 in *decode_symbol*( ) identify the symbol for which

$$cum\_freq[symbol]$$

$$\leq \left\lfloor \frac{(value - low + 1) * cum\_freq[0] - 1}{high - low + 1} \right\rfloor$$

$$< cum\_freq[symbol - 1],$$

where $\lfloor \rfloor$ denotes the "integer part of" function that comes from integer division with truncation. It is shown in the Appendix that this implies

$$low + \left\lfloor \frac{(high - low + 1) * cum\_freq[symbol]}{cum\_freq[0]} \right\rfloor$$

$$\leq v \leq low$$

$$+ \left\lfloor \frac{(high - low + 1) * cum\_freq[symbol - 1]}{cum\_freq[0]} \right\rfloor - 1,$$

so that *value* lies within the new interval that *decode_symbol*( ) calculates in lines 190–193. This is sufficient to guarantee that the decoding operation identifies each symbol correctly.

## Underflow

As Figure 1 shows, arithmetic coding works by scaling the cumulative probabilities given by the model into the interval [*low, high*] for each character transmitted. Suppose *low* and *high* are very close together—so close that this scaling operation maps some different symbols of the model onto the same integer in the [*low, high*] interval. This would be disastrous, because if such a symbol actually occurred it would not be possible to continue encoding. Consequently, the encoder must guarantee that the interval [*low, high*] is always large enough to prevent this. The simplest way to do this is to ensure that this interval is at least as large as *Max_frequency*, the maximum allowed cumulative frequency count (line 36).

How could this condition be violated? The bit-shifting operation explained above ensures that *low* and *high* can only become close together when they straddle *Half*. Suppose in fact they become as close as

$$First\_qtr \leq low < Half \leq high < Third\_qtr.$$

Then the next two bits sent will have opposite polarity, either 01 or 10. For example, if the next bit turns out to be zero (i.e., *high* descends below *Half* and [0, *Half*] is expanded to the full interval), the bit after that will be one, since the range has to be above the midpoint of the expanded interval. Conversely, if the next bit happens to be one, the one after that will be zero. Therefore the interval can safely be expanded right now, if only we remember

that, whatever bit actually comes next, its opposite must be transmitted afterwards as well. Thus lines 104–109 expand [*First_qtr, Third_qtr*] into the whole interval, remembering in *bits_to_follow* that the bit that is output next must be followed by an opposite bit. This explains why all output is done via *bit_plus_follow*( ) (lines 128–135), instead of directly with *output_bit*( ).

But what if, after this operation, it is *still* true that

$$First\_qtr \leq low < Half \leq high < Third\_qtr?$$

Figure 5 illustrates this situation, where the current [*low, high*] range (shown as a thick line) has been expanded a total of three times. Suppose the next bit turns out to be zero, as indicated by the arrow in Figure 5a being below the halfway point. Then the next three bits will be ones, since the arrow is not only in the top half of the bottom half of the original range, but in the top quarter, and moreover the top eighth, of that half—this is why the expansion can occur three times. Similarly, as Figure 5b shows, if the next bit turns out to be a one, it will be followed by three zeros. Consequently, we need only count the number of expansions and follow the next bit by that number of opposites (lines 106 and 131–134).

Using this technique the encoder can guarantee that, after the shifting operations, either

$$low < First\_qtr < Half \leq high \qquad \text{(1a)}$$

or

$$low < Half < Third\_qtr \leq high. \qquad \text{(1b)}$$

Therefore, as long as the integer range spanned by the cumulative frequencies fits into a quarter of that provided by *code_values*, the underflow problem cannot occur. This corresponds to the condition

$$Max\_frequency \leq \frac{Top\_value + 1}{4} + 1,$$

which is satisfied by Figure 3, since *Max_frequency* $= 2^{14} - 1$ and *Top_value* $= 2^{16} - 1$ (lines 36, 9). More than 14 bits cannot be used to represent cumulative frequency counts without increasing the number of bits allocated to *code_values*.

We have discussed underflow in the encoder only. Since the decoder's job, once each symbol has been decoded, is to track the operation of the encoder, underflow will be avoided if it performs the same expansion operation under the same conditions.

## Overflow

Now consider the possibility of overflow in the integer multiplications corresponding to those of Figure 2, which occur in lines 91–94 and 190–193
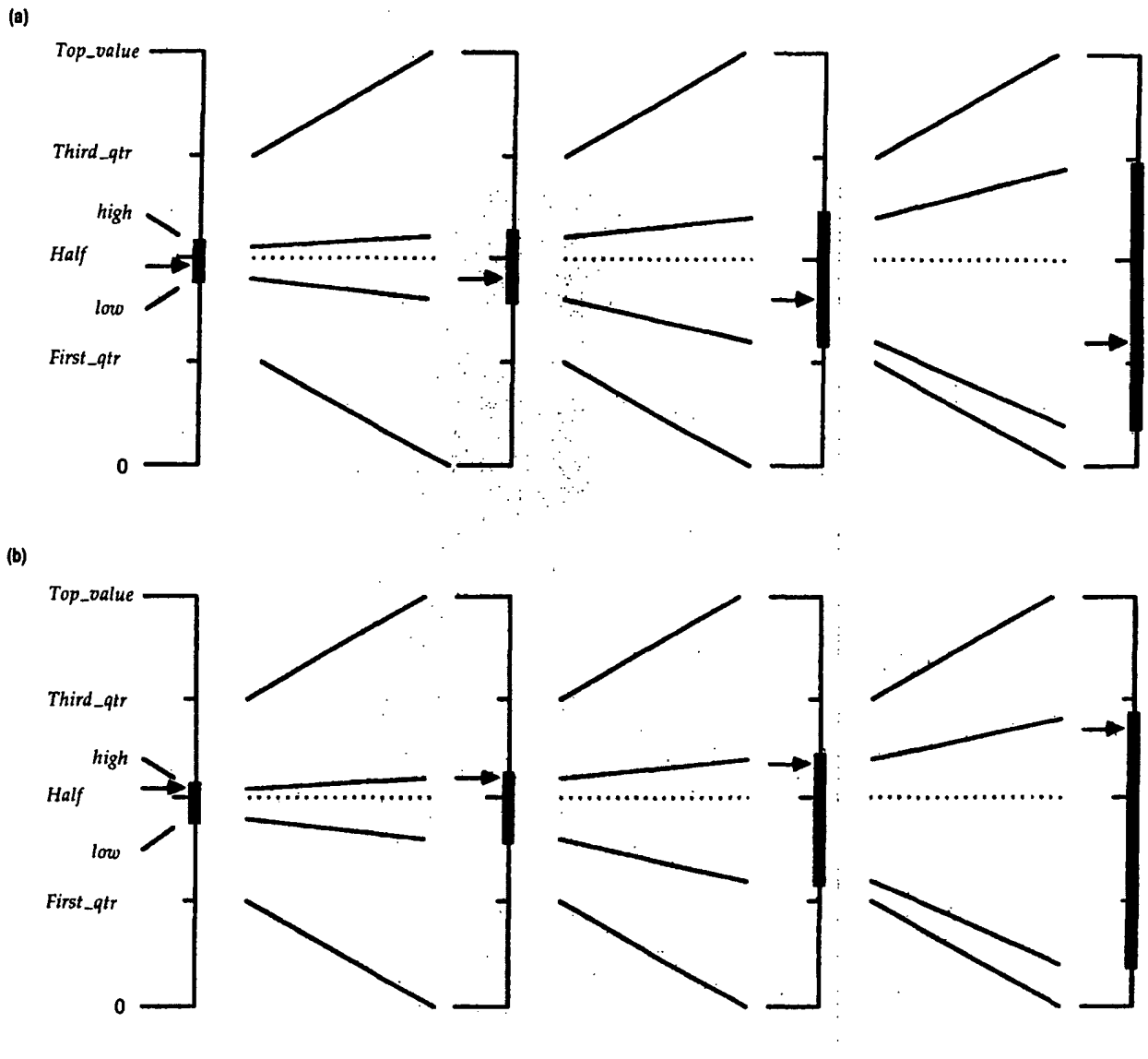
**(a)**



**(b)**

FIGURE 5. Scaling the Interval to Prevent Underflow

of Figure 3. Overflow cannot occur provided the product

$$range * Max\_frequency$$

fits within the integer word length available, since cumulative frequencies cannot exceed $Max\_frequency$. Range might be as large as $Top\_value$ + 1, so the largest possible product in Figure 3 is $2^{16}(2^{14} - 1)$, which is less than $2^{30}$. Long declarations are used for code_value (line 7) and range (lines 89, 183) to ensure that arithmetic is done to 32-bit precision.

## Constraints on the Implementation

The constraints on word length imposed by underflow and overflow can be simplified by assuming that frequency counts are represented in $f$ bits, and code_values in $c$ bits. The implementation will work correctly provided

$$f \leq c - 2$$
$$f + c \leq p,$$

the precision to which arithmetic is performed.

In most C implementations, $p = 31$ if long integers

are used, and $p = 32$ if they are *unsigned long*. In Figure 3, $f = 14$ and $c = 16$. With appropriately modified declarations, *unsigned long* arithmetic with $f = 15$ and $c = 17$ could be used. In assembly language, $c = 16$ is a natural choice because it expedites some comparisons and bit manipulations (e.g., those of lines 95–113 and 194–213).

If $p$ is restricted to 16 bits, the best values possible are $c = 9$ and $f = 7$, making it impossible to encode a full alphabet of 256 symbols, as each symbol must have a count of at least one. A smaller alphabet (e.g., the 26 letters, or 4-bit nibbles) could still be handled.

### Termination
To finish the transmission, it is necessary to send a unique terminating symbol (*EOF_symbol*, line 56) and then follow it by enough bits to ensure that the encoded string falls within the final range. Since *done_encoding*() (lines 119–123) can be sure that *low* and *high* are constrained by either Eq. (1a) or (1b) above, it need only transmit 01 in the first case or 10 in the second to remove the remaining ambiguity. It is convenient to do this using the *bit_plus_follow*() procedure discussed earlier. The *input_bit*() procedure will actually read a few more bits than were sent by *output_bit*(), as it needs to keep the low end of the buffer full. It does not matter what value these bits have, since EOF is uniquely determined by the last two bits actually transmitted.

### MODELS FOR ARITHMETIC CODING
The program in Figure 3 must be used with a model that provides a pair of translation tables *index_to_char*[ ] and *char_to_index*[ ], and a cumulative frequency array *cum_freq*[ ]. The requirements on the latter are that

- $cum\_freq[i - 1] \geq cum\_freq[i]$;
- an attempt is never made to encode a symbol $i$ for which $cum\_freq[i - 1] = cum\_freq[i]$; and
- $cum\_freq[0] \leq Max\_frequency$.

Provided these conditions are satisfied, the values in the array need bear no relationship to the actual cumulative symbol frequencies in messages. Encoding and decoding will still work correctly, although encodings will occupy less space if the frequencies are accurate. (Recall our successfully encoding *eaii!* according to the model of Table I, which does not actually reflect the frequencies in the message.)

### Fixed Models
The simplest kind of model is one in which symbol frequencies are fixed. The first model in Figure 4 has symbol frequencies that approximate those of English (taken from a part of the Brown Corpus [12]).

However, bytes that did not occur in that sample have been given frequency counts of one in case they do occur in messages to be encoded (so this model will still work for binary files in which all 256 bytes occur). Frequencies have been normalized to total 8000. The initialization procedure *start_model*() simply computes a cumulative version of these frequencies (lines 48–51), having first initialized the translation tables (lines 44–47). Execution speed would be improved if these tables were used to reorder symbols and frequencies so that the most frequent came first in the *cum_freq*[ ] array. Since the model is fixed, the procedure *update_model*(), which is called from both *encode.c* and *decode.c*, is null.

An *exact* model is one where the symbol frequencies in the message are exactly as prescribed by the model. For example, the fixed model of Figure 4 is close to an exact model for the particular excerpt of the Brown Corpus from which it was taken. To be truly exact, however, symbols that did not occur in the excerpt would be assigned counts of zero, rather than one (sacrificing the capability of transmitting messages containing those symbols). Moreover, the frequency counts would not be scaled to a predetermined cumulative frequency, as they have been in Figure 4. The exact model can be calculated and transmitted before the message is sent. It is shown by Cleary and Witten [3] that, under quite general conditions, this will *not* give better overall compression than adaptive coding (which is described next).

### Adaptive Models
An adaptive model represents the changing symbol frequencies seen *so far* in a message. Initially all counts might be the same (reflecting no initial information), but they are updated, as each symbol is seen, to approximate the observed frequencies. Provided both encoder and decoder use the same initial values (e.g., equal counts) and the same updating algorithm, their models will remain in step. The encoder receives the next symbol, encodes it, and updates its model. The decoder identifies it according to its current model and then updates its model.

The second half of Figure 4 shows such an adaptive model. This is the type of model recommended for use with Figure 3, for in practice it will outperform a fixed model in terms of compression efficiency. Initialization is the same as for the fixed model, except that all frequencies are set to one. The procedure *update_model* (*symbol*) is called by both *encode_symbol*() and *decode_symbol*() (Figure 3, lines 54 and 151) after each symbol is processed.

Updating the model is quite expensive because of the need to maintain cumulative totals. In the code

of Figure 4, frequency counts, which must be maintained anyway, are used to optimize access by keeping the array in frequency order—an effective kind of self-organizing linear search [9]. *Update_model*() first checks to see if the new model will exceed the cumulative-frequency limit, and if so scales all frequencies down by a factor of two (taking care to ensure that no count scales to zero) and recomputes cumulative values (Figure 4, lines 29–37). Then, if necessary, *update_model*() reorders the symbols to place the current one in its correct rank in the frequency ordering, altering the translation tables to reflect the change. Finally, it increments the appropriate frequency count and adjusts cumulative frequencies accordingly.

## PERFORMANCE

Now consider the performance of the algorithm of Figure 3, both in compression efficiency and execution time.

### Compression Efficiency

In principle, when a message is coded using arithmetic coding, the number of bits in the encoded string is the same as the entropy of that message with respect to the model used for coding. Three factors cause performance to be worse than this in practice:

(1)  message termination overhead;
(2)  the use of fixed-length rather than infinite-precision arithmetic; and
(3)  scaling of counts so that their total is at most *Max_frequency*.

None of these effects is significant, as we now show. In order to isolate the effect of arithmetic coding, the model will be considered to be exact (as defined above).

Arithmetic coding must send extra bits at the end of each message, causing a message termination overhead. Two bits are needed, sent by *done_encoding*() (Figure 3, lines 119–123), in order to disambiguate the final symbol. In cases where a bit stream must be blocked into 8-bit characters before encoding, it will be necessary to round out to the end of a block. Combining these, an extra 9 bits may be required.

The overhead of using fixed-length arithmetic occurs because remainders are truncated on division. It can be assessed by comparing the algorithm's performance with the figure obtained from a theoretical entropy calculation that derives its frequencies from counts scaled exactly as for coding. It is completely negligible—on the order of $10^{-4}$ bits/symbol.

The penalty paid by scaling counts is somewhat larger, but still very small. For short messages (less than $2^{14}$ bytes), no scaling need be done. Even with messages of $10^5$–$10^6$ bytes, the overhead was found experimentally to be less than 0.25 percent of the encoded string.

The adaptive model in Figure 4 scales down all counts whenever the total threatens to exceed *Max_frequency*. This has the effect of weighting recent events more heavily than events from earlier in the message. The statistics thus tend to track changes in the input sequence, which can be very beneficial. (We have encountered cases where limiting counts to 6 or 7 bits gives better results than working to higher precision.) Of course, this depends on the source being modeled. Bentley et al. [2] consider other, more explicit, ways of incorporating a recency effect.

### Execution Time

The program in Figure 3 has been written for clarity rather than for execution speed. In fact, with the adaptive model in Figure 4, it takes about 420 $\mu$s per input byte on a VAX-11/780 to encode a text file, and about the same for decoding. However, easily avoidable overheads such as procedure calls account for much of this, and some simple optimizations increase speed by a factor of two. The following alterations were made to the C version shown:

(1)  The procedures *input_bit*(), *output_bit*(), and *bit_plus_follow*() were converted to macros to eliminate procedure-call overhead.
(2)  Frequently used quantities were put in register variables.
(3)  Multiplies by two were replaced by additions (C "+=").
(4)  Array indexing was replaced by pointer manipulation in the loops at line 189 in Figure 3 and lines 49–52 of the adaptive model in Figure 4.

This mildly optimized C implementation has an execution time of 214 $\mu$s/252 $\mu$s per input byte, for encoding/decoding 100,000 bytes of English text on a VAX-11/780, as shown in Table II. Also given are corresponding figures for the same program on an Apple Macintosh and a SUN-3/75. As can be seen, coding a C source program of the same length took slightly longer in all cases, and a binary object program longer still. The reason for this will be discussed shortly. Two artificial test files were included to allow readers to replicate the results. "Alphabet" consists of enough copies of the 26-letter alphabet to fill out 100,000 characters (ending with a partially completed alphabet). "Skew statistics" contains

**TABLE II. Results for Encoding and Decoding 100,000-Byte Files**

| | Output (bytes) | VAX-11/780 | | Macintosh 512 K | | SUN-3/75 | |
|---|---|---|---|---|---|---|---|
| | | Encode time ($\mu$s) | Decode time ($\mu$s) | Encode time ($\mu$s) | Decode time ($\mu$s) | Encode time ($\mu$s) | Decode time ($\mu$s) |
| Mildly optimized C implementation | | | | | | | |
| Text file | 57,718 | 214 | 262 | 687 | 881 | 98 | 121 |
| C program | 62,991 | 230 | 288 | 729 | 950 | 105 | 131 |
| VAX object program | 73,501 | 313 | 406 | 950 | 1,334 | 145 | 190 |
| Alphabet | 59,292 | 223 | 277 | 719 | 942 | 105 | 130 |
| Skew statistics | 12,092 | 143 | 170 | 507 | 645 | 70 | 85 |
| Carefully optimized assembly-language implementation | | | | | | | |
| Text file | 57,718 | 104 | 135 | 194 | 243 | 46 | 58 |
| C program | 62,991 | 109 | 151 | 208 | 266 | 51 | 65 |
| VAX object program | 73,501 | 158 | 241 | 280 | 402 | 75 | 107 |
| Alphabet | 59,292 | 105 | 145 | 204 | 264 | 51 | 65 |
| Skew statistics | 12,092 | 63 | 81 | 126 | 160 | 28 | 36 |

Notes: Times are measured in microseconds per byte of uncompressed data.
The VAX-11/780 had a floating-point accelerator, which reduces integer multiply and divide times.
The Macintosh uses an 8-MHz MC68000 with some memory wait states.
The SUN-3/75 uses a 16.67-MHz MC68020.
All times exclude I/O and operating-system overhead in support of I/O. VAX and SUN figures give user time from the UNIX® *time* command; on the Macintosh, I/O was explicitly directed to an array.
The 4.2BSD C compiler was used for VAX and SUN; Aztec C 1.06g for Macintosh.

10,000 copies of the string *aaaabaaaac*; it demonstrates that files may be encoded into less than one bit per character (output size of 12,092 bytes = 96,736 bits). All results quoted used the adaptive model of Figure 4.

A further factor of two can be gained by reprogramming in assembly language. A carefully optimized version of Figures 3 and 4 (adaptive model) was written in both VAX and M68000 assembly languages. Full use was made of registers, and advantage taken of the 16-bit *code_value* to expedite some crucial comparisons and make subtractions of *Half* trivial. The performance of these implementations on the test files is also shown in Table II in order to give the reader some idea of typical execution speeds.

The VAX-11/780 assembly-language timings are broken down in Table III. These figures were obtained with the UNIX profile facility and are accurate only to within perhaps 10 percent. (This mechanism constructs a histogram of program counter values at real-time clock interrupts and suffers from statistical variation as well as some systematic errors.) "Bounds calculation" refers to the initial parts of *encode_symbol*() and *decode_symbol*() (Figure 3, lines 90–94 and 190–193), which contain multiply and divide operations. "Bit shifting" is the major loop in both the encode and decode routines (lines 95–113 and 194–213). The *cum* calculation in *decode_symbol*(), which requires a multiply/divide, and the following loop to identify the next symbol (lines 187–189), is "Symbol decode." Finally, "Model

**TABLE III. Breakdown of Timings for the VAX-11/780 Assembly-Language Version**

| | Encode time ($\mu$s) | Decode time ($\mu$s) |
|---|---|---|
| Text file | | |
| Bounds calculation | 32 | 31 |
| Bit shifting | 39 | 30 |
| Model update | 29 | 29 |
| Symbol decode | — | 45 |
| Other | 4 | 0 |
| | 104 | 135 |
| C program | | |
| Bounds calculation | 30 | 28 |
| Bit shifting | 42 | 35 |
| Model update | 33 | 36 |
| Symbol decode | — | 51 |
| Other | 4 | 1 |
| | 109 | 151 |
| VAX object program | | |
| Bounds calculation | 34 | 31 |
| Bit shifting | 46 | 40 |
| Model update | 75 | 75 |
| Symbol decode | — | 94 |
| Other | 3 | 1 |
| | 158 | 241 |

update" refers to the adaptive *update_model*() procedure of Figure 4 (lines 26–53).

As expected, the bounds calculation and model update take the same time for both encoding and decoding, within experimental error. Bit shifting was quicker for the text file than for the C program and object file because compression performance was

better. The extra time for decoding over encoding is due entirely to the symbol decode step. This takes longer in the C program and object file tests because the loop of line 189 was executed more often (on average 9 times, 13 times, and 35 times, respectively). This also affects the model update time because it is the number of cumulative counts that must be incremented in Figure 4, lines 49–52. In the worst case, when the symbol frequencies are uniformly distributed, these loops are executed an average of 128 times. Worst-case performance would be improved by using a more complex tree representation for frequencies, but this would likely be slower for text files.

## SOME APPLICATIONS

Applications of arithmetic coding are legion. By liberating *coding* with respect to a model from the *modeling* required for prediction, it encourages a whole new view of data compression [19]. This separation of function costs nothing in compression performance, since arithmetic coding is (practically) optimal with respect to the entropy of the model. Here we intend to do no more than suggest the scope of this view by briefly considering

(1) adaptive text compression,
(2) nonadaptive coding,
(3) compressing black/white images, and
(4) coding arbitrarily distributed integers.

Of course, as noted earlier, greater coding efficiencies could easily be achieved with more sophisticated models. Modeling, however, is an extensive topic in its own right and is beyond the scope of this article.

*Adaptive text compression* using single-character adaptive frequencies shows off arithmetic coding to good effect. The results obtained using the program in Figures 3 and 4 vary from 4.8–5.3 bits/char-

acter for short English text files ($10^3$–$10^4$ bytes) to 4.5–4.7 bits/character for long ones ($10^5$–$10^6$ bytes). Although adaptive Huffman techniques do exist (e.g., [5, 7]), they lack the conceptual simplicity of arithmetic coding. Although competitive in compression efficiency for many files, they are slower. For example, Table IV compares the performance of the mildly optimized C implementation of arithmetic coding with that of the UNIX *compact* program that implements adaptive Huffman coding using a similar model. (*Compact*'s model is essentially the same for long files, like those of Table IV, but is better for short files than the model used as an example in this article.) Casual examination of *compact* indicates that the care taken in optimization is roughly comparable for both systems, yet arithmetic coding halves execution time. Compression performance is somewhat better with arithmetic coding on all the example files. The difference would be accentuated with more sophisticated models that predict symbols with probabilities approaching one under certain circumstances (e.g., the letter $u$ following $q$).

*Nonadaptive coding* can be performed arithmetically using fixed, prespecified models like that in the first part of Figure 4. Compression performance will be better than Huffman coding. In order to minimize execution time, the total frequency count, $cum\_freq[0]$, should be chosen as a power of two so the divisions in the bounds calculations (Figure 3, lines 91–94 and 190–193) can be done as shifts. Encode/decode times of around 60 $\mu s$/90 $\mu s$ should then be possible for an assembly-language implementation on a VAX-11/780. A carefully written implementation of Huffman coding, using table lookup for encoding and decoding, would be a bit faster in this application.

*Compressing black/white images* using arithmetic coding has been investigated by Langdon and Rissanen [14], who achieved excellent results using

### TABLE IV. Comparison of Arithmetic and Adaptive Huffman Coding

| | Arithmetic coding | | | Adaptive Huffman coding | | |
|---|---|---|---|---|---|---|
| | Output (bytes) | Encode time ($\mu s$) | Decode time ($\mu s$) | Output (bytes) | Encode time ($\mu s$) | Decode time ($\mu s$) |
| Text file | 57,718 | 214 | 262 | 57,781 | 550 | 414 |
| C program | 62,991 | 230 | 288 | 63,731 | 596 | 441 |
| VAX object program | 73,546 | 313 | 406 | 76,950 | 822 | 606 |
| Alphabet | 59,292 | 223 | 277 | 60,127 | 598 | 411 |
| Skew statistics | 12,092 | 143 | 170 | 16,257 | 215 | 132 |

Notes: The mildly optimized C implementation was used for arithmetic coding.
UNIX *compact* was used for adaptive Huffman coding.
Times are for a VAX-11/780 and exclude I/O and operating-system overhead in support of I/O.

a model that conditioned the probability of a pixel's being black on a template of pixels surrounding it. The template contained a total of 10 pixels, selected from those above and to the left of the current one so that they precede it in the raster scan. This creates 1024 different possible contexts, and for each the probability of the pixel being black was estimated adaptively as the picture was transmitted. Each pixel's polarity was then coded arithmetically according to this probability. A 20–30 percent improvement in compression was attained over earlier methods. To increase coding speed, Langdon and Rissanen used an approximate method of arithmetic coding that avoided multiplication by representing probabilities as integer powers of ½. Huffman coding cannot be directly used in this application, as it never compresses with a two-symbol alphabet. Run-length coding, a popular method for use with two-valued alphabets, provides another opportunity for arithmetic coding. The model reduces the data to a sequence of lengths of runs of the same symbol (e.g., for picture coding, run-lengths of black followed by white followed by black followed by white . . .). The sequence of lengths must be transmitted. The CCITT facsimile coding standard [11] bases a Huffman code on the frequencies with which black and white runs of different lengths occur in sample documents. A fixed arithmetic code using these same frequencies would give better performance; adapting the frequencies to each particular document would be better still.

*Coding arbitrarily distributed integers* is often called for in use with more sophisticated models of text, image, or other data. Consider, for instance, the locally adaptive data compression scheme of Bentley et al. [2], in which the encoder and decoder cache the last $N$ different words seen. A word present in the cache is transmitted by sending the integer cache index. Words not in the cache are transmitted by sending a new-word marker followed by the characters of the word. This is an excellent model for text in which words are used frequently over short intervals and then fall into long periods of disuse. Their paper discusses several variable-length codings for the integers used as cache indexes. Arithmetic coding allows *any* probability distribution to be used as the basis for a variable-length encoding, including—among countless others—the ones implied by the particular codes discussed there. It also permits use of an adaptive model for cache indexes, which is desirable if the distribution of cache hits is difficult to predict in advance. Furthermore, with arithmetic coding, the code spaced allotted to the cache indexes can be scaled down to accommo-

date any desired probability for the new-word marker.

## APPENDIX. Proof of Decoding Inequality

Using one-letter abbreviations for *cum_freq, symbol, low, high*, and *value*, suppose

$$c[s] \leq \left\lfloor \frac{(v - l + 1) \times c[0] - 1}{h - l + 1} \right\rfloor$$
$$< c[s - 1];$$

in other words,

$$c[s] \leq \frac{(v - l + 1) \times c[0] - 1}{r} - e \qquad (1)$$
$$\leq c[s - 1] - 1,$$

where

$$r = h - l + 1, \qquad 0 \leq e \leq \frac{r - 1}{r}.$$

(The last inequality of Eq. (1) derives from the fact that $c[s - 1]$ must be an integer.) Then we wish to show that $l' \leq v \leq h'$, where $l'$ and $h'$ are the updated values for *low* and *high* as defined below.

(a) $\quad l' \equiv l + \left\lfloor \dfrac{r \times c[s]}{c[0]} \right\rfloor$

$$\leq l + \frac{r}{c[0]} \left[ \frac{(v - l + 1) \times c[0] - 1}{r} - e \right]$$

from Eq. (1),

$$\leq v + 1 - \frac{1}{c[0]},$$

so $l' \leq v$ since both $v$ and $l'$ are integers and $c[0] > 0$.

(b) $\quad h' \equiv l + \left\lfloor \dfrac{r \times c[s-1]}{c[0]} \right\rfloor - 1$

$$\geq l + \frac{r}{c[0]} \left[ \frac{(v - l + 1) \times c[0] - 1}{r} + 1 - e \right] - 1$$

from Eq. (1),

$$\geq v + \frac{r}{c[0]} \left[ -\frac{1}{r} + 1 - \frac{r - 1}{r} \right] = v.$$

**REFERENCES**
1. Abramson, N. *Information Theory and Coding.* McGraw-Hill, New York, 1963. This textbook contains the first reference to what was to become the method of arithmetic coding (pp. 61–62).
2. Bentley, J.L., Sleator, D.D., Tarjan, R.E., and Wei, V.K. A locally adaptive data compression scheme. *Commun. ACM 29*, 4 (Apr. 1986), 320–330. Shows how recency effects can be incorporated explicitly into a text compression system.

3. Cleary, J.C., and Witten, I.H. A comparison of enumerative and adaptive codes. *IEEE Trans. Inf. Theory IT-30*, 2 (Mar. 1984), 306–315. Demonstrates under quite general conditions that adaptive coding outperforms the method of calculating and transmitting an exact model of the message first.

4. Cleary, J.C., and Witten, I.H. Data compression using adaptive coding and partial string matching. *IEEE Trans. Commun. COM-32*, 4 (Apr. 1984), 396–402. Presents an adaptive modeling method that reduces a large sample of mixed-case English text to around 2.2 bits/character when arithmetically coded.

5. Cormack, G.V., and Horspool, R.N. Algorithms for adaptive Huffman codes. *Inf. Process. Lett. 18*, 3 (Mar. 1984), 159–166. Describes how adaptive Huffman coding can be implemented efficiently.

6. Cormack, G.V., and Horspool, R.N. Data compression using dynamic Markov modeling. Res. Rep., Computer Science Dept., Univ. of Waterloo, Ontario, Apr. 1985. Also to be published in *Comput. J.* Presents an adaptive state-modeling technique that, in conjunction with arithmetic coding, produces results competitive with those of [4].

7. Gallager, R.G. Variations on a theme by Huffman. *IEEE Trans. Inf. Theory IT-24*, 6 (Nov. 1978), 668–674. Presents an adaptive Huffman coding algorithm, and derives new bounds on the redundancy of Huffman codes.

8. Held, G. *Data Compression: Techniques and Applications.* Wiley, New York, 1984. Explains a number of ad hoc techniques for compressing text.

9. Hester, J.H., and Hirschberg, D.S. Self-organizing linear search. *ACM Comput. Surv. 17*, 3 (Sept. 1985), 295–311. A general analysis of the technique used in the present article to expedite access to an array of dynamically changing frequency counts.

10. Huffman, D.A. A method for the construction of minimum-redundancy codes. *Proc. Inst. Electr. Radio Eng. 40*, 9 (Sept. 1952), 1098–1101. The classic paper in which Huffman introduced his famous coding method.

11. Hunter, R., and Robinson, A.H. International digital facsimile coding standards. *Proc. Inst. Electr. Electron. Eng. 68*, 7 (July 1980), 854–867. Describes the use of Huffman coding to compress run lengths in black/white images.

12. Kucera, H., and Francis, W.N. *Computational Analysis of Present-Day American English.* Brown University Press, Providence, R.I., 1967. This large corpus of English is often used for computer experiments with text, including some of the evaluation reported in the present article.

13. Langdon, G.G. An introduction to arithmetic coding. *IBM J. Res. Dev. 28*, 2 (Mar. 1984), 135–149. Introduction to arithmetic coding from the point of view of hardware implementation.

14. Langdon, G.G., and Rissanen, J. Compression of black-white images with arithmetic coding. *IEEE Trans. Commun. COM 29*, 6 (June 1981), 858–867. Uses a modeling method specially tailored to black/white pictures, in conjunction with arithmetic coding, to achieve excellent compression results.

15. Pasco, R. Source coding algorithms for fast data compression. Ph.D. thesis, Dept. of Electrical Engineering, Stanford Univ., Stanford, Calif., 1976. An early exposition of the idea of arithmetic coding, but lacking the idea of incremental operation.

16. Rissanen, J.J. Generalized Kraft inequality and arithmetic coding. *IBM J. Res. Dev. 20* (May 1976), 198–203. Another early exposition of the idea of arithmetic coding.

17. Rissanen, J.J. Arithmetic codings as number representations. *Acta Polytech. Scand. Math. 31* (Dec. 1979), 44–51. Further develops arithmetic coding as a practical technique for data representation.

18. Rissanen, J., and Langdon, G.G. Arithmetic coding. *IBM J. Res. Dev. 23*, 2 (Mar. 1979), 149–162. Describes a broad class of arithmetic codes.

19. Rissanen, J., and Langdon, G.G. Universal modeling and coding. *IEEE Trans. Inf. Theory IT-27*, 1 (Jan. 1981), 12–23. Shows how data compression can be separated into *modeling* for prediction and *coding* with respect to a model.

20. Rubin, F. Arithmetic stream coding using fixed precision registers. *IEEE Trans. Inf. Theory IT-25*, 6 (Nov. 1979), 672–675. One of the first papers to present all the essential elements of practical arithmetic coding, including fixed-point computation and incremental operation.

21. Shannon, C.E., and Weaver, W. *The Mathematical Theory of Communication.* University of Illinois Press, Urbana, Ill., 1949. A classic book that develops communication theory from the ground up.

22. Welch, T.A. A technique for high-performance data compression. *Computer 17*, 6 (June 1984), 8–19. A very fast coding technique based on the method of [23], but whose compression performance is poor by the standards of [4] and [6]. An improved implementation of this method is widely used in UNIX systems under the name *compress.*

23. Ziv, J., and Lempel, A. Compression of individual sequences via variable-rate coding. *IEEE Trans. Inf. Theory IT-24*, 5 (Sept. 1978), 530–536. Describes a method of text compression that works by replacing a substring with a pointer to an earlier occurrence of the same substring. Although it performs quite well, it does not provide a clear separation between modeling and coding.

Authors' Present Address: Ian H. Witten, Radford M. Neal, and John G. Cleary, Dept. of Computer Science, The University of Calgary, 2500 University Drive NW, Calgary, Canada T2N 1N4.

---

### In response to membership requests . . .

## CURRICULA RECOMMENDATIONS FOR COMPUTING

Volume I: Curricula Recommendations for Computer Science

Volume II: Curricula Recommendations for Information Systems

Volume III: Curricula Recommendations for Related Computer Science Programs in Vocational-Technical Schools, Community and Junior Colleges and Health Computing

Information available from Deborah Cotton—Single Copy Sales (212) 869-7440 ext. 309